

# Process (ETL Project)

- Created Python virtual environment
  - Inside the project folder
- Activate venv
  - By default it is already activated (VSCode)
- pip install apache-airflow
- Created kafka\_stream.py file
  - Call api and fetch the data with a proper format
- Create docker-compose.yml file
  - To keep everything in the Docker container
  - It will keep checks on all our installations so we can run it on other machines

## 1. Install Dependencies

Before installing services, ensure you have the necessary system packages:

- `sudo apt update && sudo apt upgrade -y`
- `sudo apt install -y openjdk-11-jdk python3 python3-pip wget curl netcat`

## 2. Install and Configure Services Manually

You'll need to install and start the following services:

### A. Install and Start Zookeeper

**Install:**

- `sudo apt install -y zookeeper zookeeperd`

**Start Zookeeper:**

- `sudo systemctl start zookeeper`
- `sudo systemctl enable zookeeper`

**Check if Zookeeper is Running:**

- `echo "ruok" | nc localhost 2181`
  - If it responds with imok, Zookeeper is running successfully.

## B. Install and Start Kafka

### Download Kafka:

- `wget https://downloads.apache.org/kafka/3.6.0/kafka_2.13-3.6.0.tgz`
- `tar -xvzf kafka_2.13-3.6.0.tgz`
- `mv kafka_2.13-3.6.0 kafka`
- `cd kafka`

### Start Kafka Server:

- `bin/kafka-server-start.sh config/server.properties`

To test Kafka, create a topic:

- `bin/kafka-topics.sh --create --topic test_topic --bootstrap-server localhost:9092`
- `bin/kafka-topics.sh --list --bootstrap-server localhost:9092`

## C. Install and Start Schema Registry

### Download Confluent Schema Registry:

- `wget https://packages.confluent.io/archive/7.4/confluent-7.4.0.tar.gz`
- `tar -xvzf confluent-7.4.0.tar.gz`
- `mv confluent-7.4.0 confluent`
- `cd confluent`

### Start Schema Registry:

- `bin/schema-registry-start etc/schema-registry/schema-registry.properties`

## D. Install Apache Airflow

- `pip3 install apache-airflow`
- `export AIRFLOW_HOME=~/.airflow`
- `airflow db init`

Start the webserver:

- `airflow webserver --port 8080`

Start the scheduler:

- `airflow scheduler`

## E. Install and Start PostgreSQL

- `sudo apt install -y postgresql postgresql-contrib`
- `sudo systemctl start postgresql`
- `sudo systemctl enable postgresql`

### Create Database and User:

- `sudo -u postgres psql`

Inside PostgreSQL shell, run:

- `CREATE DATABASE airflow;`
- `CREATE USER airflow WITH ENCRYPTED PASSWORD 'airflow';`
- `GRANT ALL PRIVILEGES ON DATABASE airflow TO airflow;`
- `\q`

## F. Install Apache Spark

### Download Spark:

- `wget https://archive.apache.org/dist/spark/spark-3.5.0/spark-3.5.0-bin-hadoop3.tgz`
- `tar -xvzf spark-3.5.0-bin-hadoop3.tgz`
- `mv spark-3.5.0-bin-hadoop3 spark`

### Start Spark Master:

- `./spark/sbin/start-master.sh`

### Start Spark Worker:

- `./spark/sbin/start-worker.sh spark://localhost:7077`

## G. Install and Start Cassandra

### Install Cassandra:

- `wget https://downloads.apache.org/cassandra/4.1.3/apache-cassandra-4.1.3-bin.tar.gz`
- `tar -xvzf apache-cassandra-4.1.3-bin.tar.gz`
- `mv apache-cassandra-4.1.3 cassandra`

### Start Cassandra:

- `./cassandra/bin/cassandra -R`

To verify:

- `./cassandra/bin/cqlsh`
  - If it enters the Cassandra shell, it's working fine.

### 3. Configure Services to Work Together

Each service must be properly configured to interact with others.

- **Kafka should connect to Zookeeper (localhost:2181)**
  - **Schema Registry should connect to Kafka (localhost:9092)**
  - **Airflow should use PostgreSQL (localhost:5432)**
  - **Spark should be able to fetch data from Kafka**
  - **Processed data should be stored in Cassandra (localhost:9042)**
- 

### 4. Run Your Pipeline

Now that everything is installed and configured:

1. Start Zookeeper
2. Start Kafka
3. Start Schema Registry
4. Start Airflow (Webserver & Scheduler)
5. Start PostgreSQL
6. Start Spark Master & Worker
7. Start Cassandra

Then run your DAGs and Kafka consumers to process and store data.

## If want to keep it all inside Docker for easy installations:

version: '3'

services:

zookeeper:

image: confluentinc/cp-zookeeper:7.4.0

hostname: zookeeper

container\_name: zookeeper

ports:

- "2181:2181"

environment:

ZOOKEEPER\_CLIENT\_PORT: 2181

ZOOKEEPER\_TICK\_TIME: 2000

healthcheck:

test: ['CMD', 'bash', '-c', "echo 'ruok' | nc localhost 2181"]

interval: 10s

timeout: 5s

retries: 5

networks:

- confluent

broker:

image: confluentinc/cp-server:7.4.0

hostname: broker

container\_name: broker

depends\_on:

zookeeper:

condition: service\_healthy

ports:

- "9092:9092"

- "9101:9101"

environment:

KAFKA\_BROKER\_ID: 1

KAFKA\_ZOOKEEPER\_CONNECT: 'zookeeper:2181'

KAFKA\_LISTENER\_SECURITY\_PROTOCOL\_MAP:

PLAINTEXT:PLAINTEXT,PLAINTEXT\_HOST:PLAINTEXT

KAFKA\_ADVERTISED\_LISTENERS:

PLAINTEXT://broker:29092,PLAINTEXT\_HOST://localhost:9092

KAFKA\_METRIC\_REPORTERS: io.confluent.metrics.reporter.ConfluentMetricsReporter

KAFKA\_OFFSETS\_TOPIC\_REPLICATION\_FACTOR: 1

KAFKA\_GROUP\_INITIAL\_REBALANCE\_DELAY\_MS: 0

KAFKA\_CONFLUENT\_LICENSE\_TOPIC\_REPLICATION\_FACTOR: 1

KAFKA\_CONFLUENT\_BALANCER\_TOPIC\_REPLICATION\_FACTOR: 1

KAFKA\_TRANSACTION\_STATE\_LOG\_MIN\_ISR: 1

KAFKA\_TRANSACTION\_STATE\_LOG\_REPLICATION\_FACTOR: 1  
KAFKA\_JMX\_PORT: 9101  
KAFKA\_JMX\_HOSTNAME: localhost  
KAFKA\_CONFLUENT\_SCHEMA\_REGISTRY\_URL: http://schema-registry:8081  
CONFLUENT\_METRICS\_REPORTER\_BOOTSTRAP\_SERVERS: broker:29092  
CONFLUENT\_METRICS\_REPORTER\_TOPIC\_REPLICAS: 1  
CONFLUENT\_METRICS\_ENABLE: 'false'  
CONFLUENT\_SUPPORT\_CUSTOMER\_ID: 'anonymous'

networks:

- confluent

healthcheck:

test: [ "CMD", "bash", "-c", 'nc -z localhost 9092' ]

interval: 10s

timeout: 5s

retries: 5

schema-registry:

image: confluentinc/cp-schema-registry:7.4.0

hostname: schema-registry

container\_name: schema-registry

depends\_on:

broker:

condition: service\_healthy

ports:

- "8081:8081"

environment:

SCHEMA\_REGISTRY\_HOST\_NAME: schema-registry

SCHEMA\_REGISTRY\_KAFKASTORE\_BOOTSTRAP\_SERVERS: 'broker:29092'

SCHEMA\_REGISTRY\_LISTENERS: http://0.0.0.0:8081

networks:

- confluent

healthcheck:

test: [ "CMD", "curl", "-f", "http://localhost:8081/" ]

interval: 30s

timeout: 10s

retries: 5

control-center:

image: confluentinc/cp-enterprise-control-center:7.4.0

hostname: control-center

container\_name: control-center

depends\_on:

broker:

condition: service\_healthy

schema-registry:

condition: service\_healthy

ports:

- "9021:9021"

environment:

CONTROL\_CENTER\_BOOTSTRAP\_SERVERS: 'broker:29092'  
CONTROL\_CENTER\_SCHEMA\_REGISTRY\_URL: "http://schema-registry:8081"  
CONTROL\_CENTER\_REPLICATION\_FACTOR: 1  
CONTROL\_CENTER\_INTERNAL\_TOPICS\_PARTITIONS: 1  
CONTROL\_CENTER\_MONITORING\_INTERCEPTOR\_TOPIC\_PARTITIONS: 1  
CONFLUENT\_METRICS\_TOPIC\_REPLICATION: 1  
CONFLUENT\_METRICS\_ENABLE: 'false'  
PORT: 9021

networks:

- confluent

healthcheck:

test: [ "CMD", "curl", "-f", "http://localhost:9021/health" ]  
interval: 30s  
timeout: 10s  
retries: 5

webserver:

image: apache/airflow:2.6.0-python3.9  
command: webserver  
entrypoint: ["/opt/airflow/script/entrypoint.sh"]  
depends\_on:  
- postgres  
environment:  
- LOAD\_EX=n  
- EXECUTOR=Sequential  
-

AIRFLOW\_\_DATABASE\_\_SQL\_ALCHEMY\_CONN=postgresql+psycopg2://airflow:airflow@postgres:5432/airflow

- AIRFLOW\_WEBSERVER\_SECRET\_KEY=this\_is\_a\_very\_secured\_key

logging:

options:  
max-size: 10m  
max-file: "3"

volumes:

- ./dags:/opt/airflow/dags
- ./script/entrypoint.sh:/opt/airflow/script/entrypoint.sh
- ./requirements.txt:/opt/airflow/requirements.txt

ports:

- "8080:8080"

healthcheck:

test: ['CMD-SHELL', "[ -f /opt/airflow/airflow-webserver.pid ]"]  
interval: 30s  
timeout: 30s  
retries: 3

networks:

- confluent

scheduler:

image: apache/airflow:2.6.0-python3.9

depends\_on:

webserver:

condition: service\_healthy

volumes:

- ./dags:/opt/airflow/dags
- ./script/entrypoint.sh:/opt/airflow/script/entrypoint.sh
- ./requirements.txt:/opt/airflow/requirements.txt

environment:

- LOAD\_EX=n
- EXECUTOR=Sequential

-

AIRFLOW\_\_DATABASE\_\_SQL\_ALCHEMY\_CONN=postgresql+psycopg2://airflow:airflow@postgres:5432/airflow

- AIRFLOW\_\_WEBSERVER\_\_SECRET\_KEY=this\_is\_a\_very\_secured\_key

command: bash -c "pip install -r ./requirements.txt && airflow db upgrade && airflow scheduler"

networks:

- confluent

postgres:

image: postgres:14.0

environment:

- POSTGRES\_USER=airflow
- POSTGRES\_PASSWORD=airflow
- POSTGRES\_DB=airflow

logging:

options:

max-size: 10m

max-file: "3"

networks:

- confluent

spark-master:

image: bitnami/spark:latest

command: bin/spark-class org.apache.spark.deploy.master.Master

ports:

- "9090:8080"
- "7077:7077"

networks:

- confluent

spark-worker:

image: bitnami/spark:latest

command: bin/spark-class org.apache.spark.deploy.worker.Worker

spark://spark-master:7077

depends\_on:

- spark-master

environment:



SPARK\_MODE: worker  
SPARK\_WORKER\_CORES: 2  
SPARK\_WORKER\_MEMORY: 1g  
SPARK\_MASTER\_URL: spark://spark-master:7077

networks:

- confluent

cassandra\_db:

image: cassandra:latest  
container\_name: cassandra  
hostname: cassandra

ports:

- "9042:9042"

environment:

- MAX\_HEAP\_SIZE=512M
- HEAP\_NEWSIZE=100M
- CASSANDRA\_USERNAME=cassandra
- CASSANDRA\_PASSWORD=cassandra

networks:

- confluent

networks:

confluent: